

WHAT IS CLAIMED IS:

1. A system comprising:
a network including a plurality of computing nodes;
5 a plurality of replicas of an object, wherein the plurality of replicas are stored
on a first plurality of the nodes;
wherein the network includes a first node operable to initiate an update
operation to update the plurality of replicas of the object, wherein said updating the
plurality of replicas of the object comprises updating a subset but not all of the replicas;
10 wherein for each node on which one of the replicas was updated in the update
operation, the node is operable to add the object to a list of incoherent objects.
2. The system of claim 1,
wherein said initiating the update operation to update the plurality of replicas
15 of the object comprises initiating an update operation to attempt to update all of the
replicas;
wherein only the subset of the replicas are updated because one or more of the
replicas are unreachable.
- 20 3. The system of claim 1,
wherein said updating the subset of the replicas includes updating a first
replica, wherein the first replica is stored on a second node;
wherein after said adding the object to the list of incoherent objects, the second
node is operable to attempt to communicate with all the replicas of the object;
25 wherein if all the replicas of the object are reachable then the replicas that were
not in the subset of replicas that were updated are synchronized with the replicas that
were updated.
4. The system of claim 3,

wherein said updating the subset of the replicas includes applying a first change to each replica in the subset of the replicas;

wherein said synchronizing the replicas that were not in the subset of replicas that were updated with the replicas that were updated comprises applying the first change
5 to each replica that was not in the subset.

5. The system of claim 4,

wherein after said synchronizing, each node in the subset of replicas that were updated is operable to remove the object from its list of incoherent objects.

10

6. The system of claim 3,

wherein the list of incoherent objects on the second node includes a plurality of objects, wherein each object has a plurality of replicas, wherein the plurality of replicas for each object are stored on a plurality of nodes;

15

wherein for each object in the list of incoherent objects on the second node, the second node is operable to attempt to communicate with all the replicas of the object to synchronize the replicas.

7. The system of claim 3,

20

wherein the second node is operable to periodically perform said attempting to communicate with all the replicas of the object.

8. The system of claim 7,

wherein if said periodically attempting to communicate with all the replicas of
25 the object has not succeeded after a first amount of time has passed, the second node is operable to initiate an operation to create one or more new replicas of the object to replace one or more unreachable replicas of the object.

9. The system of claim 1,

wherein for each node on which one of the replicas was updated in the update operation, the list of incoherent objects on the node is stored in persistent storage;

wherein for each node on which one of the replicas was updated in the update operation, said node adding the object to its list of incoherent objects does not include
5 changing the list of incoherent objects in persistent storage;

wherein for each node on which one of the replicas was updated in the update operation, the node is operable to periodically update its list of incoherent objects in persistent storage to reflect new additions to the list.

10 10. The system of claim 1,
 wherein the plurality of replicas of the object comprises a plurality of persistent replicas of the object.

 11. The system of claim 1,
15 wherein said initiating the update operation to update the plurality of replicas of the object comprises initiating a distributed transaction to update at least a quorum of the replicas;

 wherein said updating the subset of the replicas comprises updating a quorum of the replicas.

20 12. The system of claim 1,
 wherein said updating the subset of the replicas comprises applying a change to data in each replica in the subset.

25 13. A system comprising:
 a network including a plurality of computing nodes;
 a plurality of replicas of an object, wherein the plurality of replicas are stored on a first plurality of the nodes;

wherein the network includes a first node operable to store a first timestamp associated with a replica on the first node;

wherein in response to receiving a read request, the first node is operable to:

determine whether time elapsed since the time indicated by the first
5 timestamp exceeds a first threshold amount; and

if the time elapsed does not exceed the first threshold amount,
respond to the read request using the replica on the first node.

14. The system of claim 13,

10 wherein if the time elapsed does exceed the first threshold amount then the first node is operable to communicate with a second node to synchronize the replica on the first node with a replica on the second node if necessary;

wherein the first node is operable to respond to the read request using the replica on the first node after said communicating with the second node.

15

15. The system of claim 14,

wherein said synchronizing the replica on the first node with the replica on the second node comprises:

determining whether the replica on the second node has received one or more
20 updates since the replica on the first node was last known to be coherent with respect to the replica on the second node; and

applying the one or more updates to the replica on the first node if there are one or more updates.

25 16. The system of claim 15,

wherein said synchronizing the replica on the first node with the replica on the second node further comprises updating the first timestamp to indicate a current time.

17. The system of claim 16,

wherein the first timestamp indicates a time at which the replica on the first node was last known to be coherent.

18. The system of claim 16,
5 wherein said updating the first timestamp to indicate the current time comprises updating the first timestamp even if the replica on the second node has not received any updates.

19. The system of claim 13,
10 wherein in response to receiving an update from a replica on a second node, the first node is operable to update the first timestamp to indicate a current time.

20. The system of claim 19,
wherein in response to receiving the update the first node is further operable to
15 apply the update to the replica on the first node.

21. A carrier medium comprising program instructions executable to implement the method of:
20 storing a plurality of replicas of an object on a plurality of nodes;
a first node initiating an update operation to update the plurality of replicas of the object, wherein said updating the plurality of replicas of the object comprises updating a subset but not all of the replicas; and
for each node on which one of the replicas was updated in the update
25 operation, the node adding the object to a list of incoherent objects.

22. The carrier medium of claim 21,
wherein said initiating the update operation to update the plurality of replicas of the object comprises initiating an update operation to attempt to update all of the
30 replicas;

wherein only the subset of the replicas are updated because one or more of the replicas are unreachable.

23. The carrier medium of claim 21,
5 wherein said updating the subset of the replicas includes updating a first replica, wherein the first replica is stored on a second node;
wherein after said adding the object to the list of incoherent objects, the second node is operable to attempt to communicate with all the replicas of the object;
wherein if all the replicas of the object are reachable then the replicas that were
10 not in the subset of replicas that were updated are synchronized with the replicas that were updated.

24. The carrier medium of claim 23,
wherein the second node is operable to remove the object from its list of
15 incoherent objects after said synchronizing.

25. The carrier medium of claim 23,
wherein the second node is operable to periodically perform said attempting to communicate with all the replicas of the object.

20
26. The carrier medium of claim 25,
wherein if said periodically attempting to communicate with all the replicas of the object has not succeeded after a first amount of time has passed, the second node is operable to initiate an operation to create one or more new replicas of the object to
25 replace one or more unreachable replicas of the object.

27. A carrier medium comprising program instructions executable to implement the method of:

storing a plurality of replicas of an object on a plurality of nodes, wherein said storing the plurality of replicas includes storing a first replica on a first node;

the first node storing a first timestamp associated with the first replica on the first node; and

5 in response to receiving a read request, the first node:

determining whether time elapsed since the time indicated by the first timestamp exceeds a first threshold amount; and

if the time elapsed does not exceed the first threshold amount, responding to the read request using the first replica on the first node.

10

28. The carrier medium of claim 27,

wherein if the time elapsed does exceed the first threshold amount then the first node communicates with a second node to synchronize the first replica on the first node with a second replica on the second node if necessary;

15 wherein the first node responds to the read request using the first replica on the first node after said communicating with the second node.

29. The carrier medium of claim 28,

20 wherein said synchronizing the first replica on the first node with the second replica on the second node comprises:

determining whether the second replica on the second node has received one or more updates since the first replica on the first node was last known to be coherent with respect to the second replica on the second node; and

25 applying the one or more updates to the first replica on the first node if there are one or more updates.

30. The carrier medium of claim 27,

wherein in response to receiving an update from a second replica on a second node, the first node updates the first timestamp to indicate a current time.

30

31. The carrier medium of claim 30,
wherein in response to receiving the update the first node is further operable to
apply the update to the first replica on the first node.

5